

On the Throughput Capacity of Information-Centric Networks

Bitan Azimdoost^{†*}, Cedric Westphal^{†*}, and Hamid R. Sadjadpour[†]

[†]Department of Electrical Engineering and [‡]Computer Engineering
University of California Santa Cruz, Santa Cruz, CA 95064, USA
{bazimdoost,cedric,hamid}@soe.ucsc.edu

* Huawei Innovation Center, Santa Clara, CA 95050, USA
{bita.azimdoost,cedric.westphal}@huawei.com

Abstract—Information-centric networks make storage one of the network primitives, and propose to cache data within the network in order to improve latency to access content and reduce bandwidth consumption. We study the throughput capacity of an information-centric network when the data cached in each node has a limited lifetime. The results show that with some fixed request and cache expiration rates the network can have the maximum throughput order of $1/\sqrt{n}$ and $1/\log n$ in case of grid and random networks, respectively. Comparing these values with the corresponding throughput with no cache capability ($1/n$ and $1/\sqrt{n \log n}$ respectively), we can actually quantify the asymptotic advantage of caching. Moreover, since the request rates will decrease as a result of increasing download delays, increasing the content lifetimes according to the network growth may result in higher throughput capacities.

I. INTRODUCTION

In today's networking situations, users are mostly interested in accessing content regardless of which host is providing this content. They are looking for a fast and secure access to data in a whole range of situations: wired or wireless; heterogeneous technologies; in a fixed location or when moving. The dynamic characteristics of the network users makes the host-centric networking paradigm inefficient. Information-centric networking (ICN) is a new networking architecture where content is accessed based upon its name, and independently of the location of the hosts [1]–[4]. In most ICN architectures, data is allowed to be stored in the nodes and routers within the network in addition to the content publisher's servers. This reduces the burden on the servers and on the network operator, and shortens the access time to the desired content.

Combining content routing with in-network-storage for the information is intuitively attractive, but there has been few works considering the impact of such architecture on the capacity of the network in a formal or analytical manner. In this work we study an information-centric network where nodes can both route and cache content. We also assume that a node will keep a copy of the content only for a finite period of time, that is until it runs out of memory space in its cache and has to rotate content, or until it ceases to serve a specific content.

The nodes issue some queries for content that is not locally available. We suppose that there exists a server which permanently keeps all the contents. This means that the content

is always provided at least by its publisher, in addition to the potential copies distributed throughout the network. Therefore, at least one replica of each content always exists in the network and if a node requests a piece of information, this data will be furnished either by its original server or by a cache containing the desired data. When the customer receives the content, it will store the content and share it with the other nodes if needed.

The present paper thus investigates the throughput capacity in such content-centric networks and addresses the following questions:

- 1) Looking at the throughput capacity, can we quantify the performance improvement brought about by a content-centric network architecture over networks with no content sharing capability?
- 2) How does the caching policy, and in particular, the length of time each piece of content spends in the cache's memory, affect the capacity?

We state two Theorems below. Theorem 1 will answer the first question studying two different network models (grid and random network) and two content discovery scenarios (shortest path to the server and flooding). Theorem 2 derives some conditions on the respective request rate (namely, the popularity of the content) and the time spent in the cache, so that these throughputs can be supported by all the nodes and the flow in no node be a bottleneck. These theorems demonstrate that adding the content sharing capability to the nodes can significantly increase the capacity.

Theorem 1. Consider a wireless network consisting of n nodes, with each node containing the information in its local cache with probability ρ . Each node can transmit at W bits per second over a common wireless channel shared by all nodes.

- Scenario i- If the nodes are located on a grid and search for the contents just on the shortest path toward the server, the maximum achievable throughput capacity order¹ is

¹ $f(n) = O(g(n))$ if $\sup_n (f(n)/g(n)) < \infty$. $f(n) = \Omega(g(n))$ if $g(n) = O(f(n))$. $f(n) = \Theta(g(n))$ if both $f(n) = O(g(n))$ and $f(n) = \Omega(g(n))$.

$$\gamma_{max} = \begin{cases} \Theta(\frac{W\rho}{\sqrt{n(1-\rho)}}), & \rho = \Omega(n^{-1/2}) \\ \Theta(\frac{W}{n(1-\rho)}), & \rho = O(n^{-1/2}) \end{cases}$$

- *Scenario ii-* If the nodes are located on a grid and use flooding as their content search algorithm, the maximum achievable throughput is

$$\gamma_{max} = \begin{cases} \Theta(\frac{W\rho^{0.4646}}{\sqrt{n(1-\rho)}}), & \rho = \Omega(n^{-1/2}) \\ \Theta(\frac{W}{n(1-\rho)}), & \rho = O(n^{-1/2}) \end{cases}$$

- *Scenario iii-* If the nodes are randomly distributed over a unit square area and use path-wise content discovery algorithm, the maximum achievable capacity is

$$\gamma_{max} = \begin{cases} \Theta(\frac{W}{\log n(1-\rho)}), & \rho = \Omega(\log n^{-1}) \\ \Theta(\frac{W}{\sqrt{n \log n(1-\rho)}}), & \rho = O(\log n^{-1}) \end{cases}$$

Theorem 2. Assume that in networks of Theorem 1 the content request and dropping rates are λ and μ , respectively. The throughput capacities of Theorem 1 are supportable if $\frac{\lambda}{\mu} = O(\frac{n \log \log n}{\log n})$ in scenario i, ii, and $\frac{\lambda}{\mu} = O(\frac{\log n \log \log \log n}{\log \log n})$ in scenario iii.

The rest of the paper is organized as follows. After a brief review of the related work, the network models, the content availability and the content discovery algorithms used in the current work are introduced in Section II. Theorems 1 and 2 are proved in Sections III, IV, respectively. Finally we will discuss the results in Section V.

II. PRELIMINARIES

A. Related Work

Many aspects of ICN networks have been studied in prior works [3]. Some performance metrics like miss ratio in the cache, or the average number of hops each request travels to locate the content have been studied in [5], [6].

Optimal cache locations [7] and cache replacement techniques [8] are two other aspects most commonly investigated. And an analytical framework for investigating properties of these networks like fairness of cache usage is proposed in [9]. [10] considered information being cached for a limited amount of time at each node, as we do here, but focused on flooding mechanism to locate the content, not on the capacity of the network.

However, to the best of our knowledge, there are just a few works focusing on the achievable data rates in such networks. [11] uses a network simulation model and evaluates the performance (file transfer delay) in a cache-and-forward system with no request for the data. [12] proposes an analytical model for single cache miss probability and stationary throughput in cascade and binary tree topologies. [13] considers a general problem of delivering content cached in a wireless network and provides some bounds on the caching capacity region from an information-theoretic point of view. Some scaling regimes for the required link capacity is computed in [14] for a static cache placement in a multihop wireless network.

B. Network Model

Two network models are studied in this work.

1) *Grid Network:* Assume that the network consists of n nodes $V = \{v_1, v_2, \dots, v_n\}$ each with a local cache of size L located on a grid (Figure 1). The distance between two adjacent nodes equals to the transmission range of each node, so the packets sent from a node are only received by four adjacent nodes. There are m different contents, $F = \{f_1, f_2, \dots, f_m\}$ with sizes B_i , $i = 1, \dots, m$, for which each node v_j may issue a query. Based on the content discovery algorithms which will be explained later in this section, the query will be transmitted in the network to discover a node containing the desired content locally. v_j then downloads b bits of data with rate γ in a hop-by-hop manner through the path P_{xj} from either a node ($v_i, x = i$) containing it locally ($f \in v_i$) or the server ($x = s$). When the download is completed, the end user stores the data in its local cache and shares it with other nodes. $P_{js \rightarrow i}$ denotes the nodes on the path from v_j to server before reaching node v_i . Without loss of generality, we assume that the server is attached to the node located at the middle of the network, changing the location of the server does not affect the scaling laws. Using the protocol model and according to [15] the transport capacity in such network is upper bounded by $\Theta(W\sqrt{n})$. This is the model studied in the first two scenarios

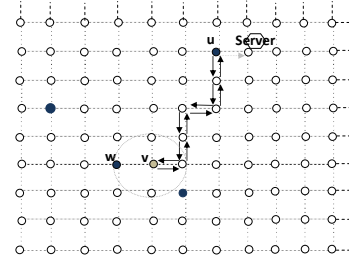


Fig. 1. The transmission range of node v contains four surrounding nodes. The black vertices contain the content in their local caches. The arrow lines demonstrate a possible discovery and receive path in scenario i, where node v downloads the required information from u . In scenario ii, v will download the data from w instead.

of Theorem 1.

2) *Random Network:* The last network studied in Theorem 1 is a more general network model where the nodes are randomly distributed over a unit square area according to a uniform distribution. We use the same model used in [15] (section 5) and divide the network area into square cells each with side-length proportional to the transmission range $r(n)$, which is selected to be at least in the order of $\sqrt{\frac{\log n}{n}}$ to guarantee the connectivity of the network [16]. According to the protocol model [15], if the cells are far enough they can transmit data at the same time with no interference; we assume that there are M^2 non-interfering groups which take turn to transmit at the corresponding time-slot in a round robin fashion. The server is assumed to be located at the middle of the network. In this model the maximum number

of simultaneous feasible transmissions will be in the order of $\frac{1}{r^2(n)}$ as each transmission consumes an area proportional to $r^2(n)$.

All the other assumptions regarding the contents, requests, and time-out durations are similar to the grid network.

C. Content Discovery Algorithm

1) *Path-wise Discovery*: To discover the location of the desired content, the request is sent through the shortest path toward the server containing the requested content. If an intermediate node has the data in its local cache, it does not forward the request toward the server anymore and the requester will start downloading from the discovered cache. Otherwise, the request will go all the way toward the server and the content is obtained from the main source. In case of the random network when a node needs a piece of information, it will send a request to its neighbors toward the server, i.e. the nodes in the same cell and one adjacent cell in the path toward the server, if any copy of the data is found it will be downloaded. If not, just one node in the adjacent cell will forward the request to the next cell toward the server.

2) *Flooding/Ring Search*: In this algorithm the request for the information is sent to all the nodes in the transmission range of the requester. If a node receiving the request contains the required data in its local cache, it notifies the requester and then downloading from the discovered cache is started. Otherwise, all the nodes that receive the request will broadcast the request to their own neighbors. This process continues until the content is discovered in a cache and the downloading follows after that.

3) *Content Distribution in Steady-State*: The time diagram of data access process in the studied network is illustrated in Figure 2. When a query for content f_i is initiated, the content is available at the requester's cache after a wait time (T_3) which is a function of the distance between the user and the data source (server or an intermediate cache), the data size, and the download speed. The expiration timer will be reset upon receiving the data and this data will be dropped after an exponentially distributed holding time (T_1) with mean $1/\mu_i$. The user may re-issue a query for that data after another exponentially distributed time (T_2) with mean $1/\lambda_i$. The solid lines in this diagram denote the portions of time that the data is available at local cache.

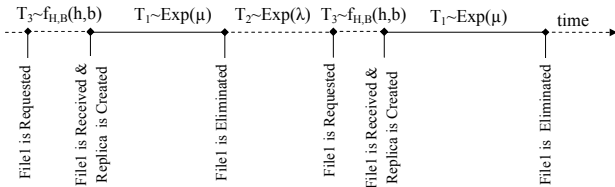


Fig. 2. Data access process time diagram in a cache network

In this work we assume identical content sizes $B_i = B$, and assume all the contents have the same popularity leading to similar request rates $\lambda_i = \lambda$, and the same time-outs $\mu_i = \mu$. As the requests for different contents are supposed to be independent and time-outs are set for each content independent of the others, we can do the calculations for one single content. If the total number of contents is not a function of the network size, this will not change the capacity order. Suppose that B is much larger than the request size, so we ignore the overhead of the discovery phase in our calculations. Furthermore, if the information sizes are the same and the download rates are also the same, the download time will be a function of the number of hops (h) between the source and the customer; $T_3 = Bh/\gamma$. In the steady-state analysis, we ignore this constant time.

We assume that each node generates a request for a content according to an exponential process with mean inter-arrival time $1/\lambda$ if it does not have it in its local cache. On the other hand, there is a time-out timer for each cached content which will reset upon receiving a content and times-out according to an exponential process with average time $1/\mu$. Therefore, each node's cache states for each piece of information are changing according to a Markov process with two states 0 and 1, and transition rates λ for change from state 0 to 1, and μ from 1 to 0.

The average portion of time that each node contains a content in its local cache is

$$\rho = \frac{1/\mu}{1/\mu + 1/\lambda} = \frac{\lambda}{\lambda + \mu}, \quad (1)$$

which is the average probability that a node contains the data at steady-state.

III. PROOF OF THEOREM 1

In this section, we prove Theorem 1 by utilizing some lemmas.

Lemma 1. Consider wireless networks described in Theorem 1. For sufficiently large networks and when ρ is large enough ($\rho = \Omega(n^{-1/2})$ for case *i*, *ii* and $\rho = \Omega(\log^{-1} n)$ for case *iii*), the average number of hops between the customer and the nearest cached content location is

$$\bar{h} = \begin{cases} \Theta(\frac{1}{\rho}), & (i) \\ \Theta(\frac{1}{\rho^{0.4646}}), & (ii) \\ \Theta(1). & (iii) \end{cases} \quad (2)$$

Proof: Scenario *i*- The probability that the first node on the path from the requester (v) to the server contains the content in its local cache is ρ , and the probability that the closest node to the customer on the query path caching the requested information is the h^{th} node is $(1 - \rho)^{h-1}\rho$. Thus the average number of hops between the customer and the

nearest cached content location is

$$\begin{aligned}\bar{h} &= \sum_{h=1}^{\sqrt{n}} hP(H=h) \\ &= \sum_{h=1}^{\sqrt{n}} h\rho(1-\rho)^{h-1} \\ &\stackrel{\text{Large } n}{\cong} \Theta\left(\frac{1}{\rho}\right)\end{aligned}\quad (3)$$

The result is valid for $\rho = \Omega(n^{-1/2})$, and for lesser values of ρ we will have $\bar{h} = \Theta(\sqrt{n})$.

Scenario *ii* - The probability that the discovered cache is located at a distance of one hop from the requester is the probability that one of the nodes on the ring at one hop distance contains the data (it consists of 4 nodes), which equals to $1 - (1 - \rho)^4$, and the probability that the data needs to travel through h hops from the discovered cache to where it is required is $(1 - (1 - \rho)^{4h}) \prod_{k=1}^{h-1} (1 - \rho)^{4k}$ as there are $4h$ nodes at distance of h hops. Therefore,

$$\begin{aligned}\bar{h} &= \sum_{h=1}^{\sqrt{n}} h(1 - (1 - \rho)^{4h}) \prod_{k=1}^{h-1} (1 - \rho)^{4k} \\ &\stackrel{\text{Large } n}{\cong} \Theta\left(\frac{1}{\rho^{0.4646}}\right)\end{aligned}\quad (4)$$

where the last equality is correct when $\rho = \Omega(n^{-1/2})$. For smaller ρ 's, \bar{h} will increase to $\Theta(\sqrt{n})$.

Scenario *iii* - The discovered cache is one hop away from the requester if there is a replica of the data in a cache at the same cell or at the adjacent cell toward the server. So since there are $\log n$ nodes in each cell, the probability of the discovered cache being at one hop distance is $1 - (1 - \rho)^{2 \log n}$, and the probability of the discovered cache being at distance of h hops away from the requester is $(1 - \rho)^{h \log n} (1 - (1 - \rho)^{\log n})$. The maximum number of hops that may be traveled this way is $\frac{1}{r(n)}$. Thus

$$\begin{aligned}\bar{h} &= 1 - (1 - \rho)^{2 \log n} \\ &+ \sum_{h=2}^{\frac{1}{r(n)}} h(1 - \rho)^{h \log n} (1 - (1 - \rho)^{\log n}) \\ &\stackrel{\text{Large } n}{\cong} \Theta(1)\end{aligned}\quad (5)$$

where the last equality is correct when $\rho = \Omega(\log^{-1} n)$. ■

In scenario *iii* the average number of hops between the nearest content location and the customer is just $\Theta(1)$ hop. This is the result of having $\log(n)$ caches in one hop distance for every requester. Each one of these caches can be a potential source for the content. When the network grows, this number will increase and if ρ is large enough ($\rho = \Omega(\log^{-1} n)$) the probability that at least one of these nodes contain the required data will approach 1, i.e., $\lim_{n \rightarrow \infty} (1 - (1 - \rho)^{\log n}) = 1$.

Lemma 2. The average probability that the server needs to

serve a request is

$$p_s = \begin{cases} \Theta\left(\frac{(2-\rho)^2}{n\rho^2}\right), & (i) \\ O\left(\frac{(2-\rho)^2}{n\rho^2}\right), & (ii) \\ \Theta\left(\frac{(2-\rho) \log n}{n}\right). & (iii) \end{cases}\quad (6)$$

Proof: Scenario *i*- The data will need to be downloaded from the server (at average distance \bar{h}_s) if no copy of the data is available on the path between a requester node and the server. As the network area is assumed to be a square and the server is in the middle of it, this probability is bounded by

$$\frac{1 + \sum_{k=1}^{h_{max}/2} 4k(1-\rho)^k}{n} \leq p_s \leq \frac{1 + \sum_{k=1}^{h_{max}} 4k(1-\rho)^k}{n}$$

Thus for large n , $p_s = \Theta\left(\frac{(2-\rho)^2}{n\rho^2}\right)$.

Note that both \bar{h}_s and h_{max} , the maximum number of hops which may be traveled between the requester and the node that possesses a valid copy of data, in this scenario are $\Theta(\sqrt{n})$.

Scenario *ii*- The data will need to be downloaded from the server (at average distance \bar{h}_s) if no copy of the content is available in the network caches. Since comparing to scenario *i* more nodes will be involved in the process of content discovery, it is obvious that in this case the request will be forwarded to the server with less probability. Thus $p_s = O\left(\frac{(2-\rho)^2}{n\rho^2}\right)$.

Scenario *iii*- The data is downloaded from the server if no node in the cells on the path toward the server cell contains a copy of the content.

$$\begin{aligned}p_s &= \frac{1 + 5 \log n (1 - \rho) + \sum_{h=2}^{\frac{1}{r(n)}} 4h \log n (1 - \rho)^{(h-1) \log n}}{n} \\ &\stackrel{\text{Large } n}{\cong} \frac{(2-\rho) \log n}{n}\end{aligned}\quad (7)$$

It can be seen that in all cases the average number of hops between the server and the node requesting the content is a function of the total number of nodes in the network and ρ .

Here we can prove *Theorem 1* using the above lemmas.

Proof: Assume that each content is retrieved with rate γ bits/sec. The traffic generated because of one download from a cache at average distance of \bar{h} hops from the requester node is $\gamma \bar{h}$ and the traffic generated due to the downloads from the server at average distance of \bar{h}_s hops from the requester is $\gamma \bar{h}_s$. The probability that the server is uploading the data is p_s and the probability that a cache node is serving the customer is $p = 1 - p_s$. The total number of requests for a content in the network at any given time is limited by the number of nodes not having the content in their own cache $((1 - \rho)n)$. Thus the maximum total bandwidth needed to accomplish these downloads will be $(1 - \rho)n(p\bar{h} + p_s\bar{h}_s)\gamma$, which is upper limited by $(\Theta(W\sqrt{n}))$ in scenarios *i*, *ii*, and $(\Theta(\frac{W}{r^2(n)}))$ in scenario *iii*.

Therefore the maximum download rate is

$$\gamma_{max} = \begin{cases} \Theta\left(\frac{W\sqrt{n}/n(1-\rho)}{(1 - \frac{(2-\rho)^2}{n\rho^2})\rho^{-1} + \frac{(2-\rho)^2}{n\rho^2}\sqrt{n}}\right), & (i) \\ \Theta\left(\frac{W\sqrt{n}/n(1-\rho)}{(1 - \frac{(2-\rho)^2}{n\rho^2})\rho^{-0.4646} + \frac{(2-\rho)^2}{n\rho^2}\sqrt{n}}\right), & (ii) \\ \Theta\left(\frac{W/r^2(n)n(1-\rho)}{(1 - \frac{(2-\rho) \log n}{n}) + \frac{(2-\rho) \log n}{nr(n)}}\right), & (iii) \end{cases}\quad (8)$$

The results of Theorem 1 can be derived by approximating these equations for sufficiently large n . Note that if there were no cache in the system, or ρ is less than the stated threshold values, all the requests would be served by the server, and the maximum download rate would be $\frac{W}{\sqrt{nh_s}} = \Theta(\frac{W}{n})$ in case i , ii and $\Theta(\frac{W}{\sqrt{n \log n}})$ in case iii . ■

IV. PROOF OF THEOREM 2

In the previous section the maximum throughput capacity in a cache wireless network has been calculated. Now it is important to verify if this throughput can be supported by each cell (node), i.e. the traffic carried by each cell (node) is not more than what it can support ($\Theta(1)$). Here we start with scenario iii and a complete proof. Scenario i will be then briefly studied. Similar reasoning can be used for scenario ii .

Proof: Scenario iii - The traffic load at the server is $\gamma_{max} p_s n(1 - \rho) = \Theta(1)$. So the flow at the server will not be a bottleneck.

The traffic load at a node as a customer will not be a bottleneck either as it does not exceed the maximum data rate which is $\gamma_{max} = \Theta(\frac{W}{\log n(1-\rho)}) < \Theta(1)$.

To compute the traffic load at a node which is serving a request, we need to know how many requests that node may serve at a time. A node $v_i \in V$ is the download source if it has the information (ρ), it is in the same cell as the requester or in a cell on the path from the requester to the server and no node in the previous cells on this path contains the content ($(1 - \rho)^{x \log n}$ where x is the number of hops between v_j and v_i), and among those nodes in the same cell which have the data v_i is selected to serve the query ($\sum_{k=1}^{\log n} \frac{1}{k} \binom{\log n - 1}{k-1} \rho^{k-1} (1 - \rho)^{\log n - k}$). For not too small ρ and large n , we have

$$P(v_i \text{ is serving } v_j's \text{ request}) = \begin{cases} \frac{1}{\log n}, & v_j \text{ and } v_i \text{ are in the same cell} \\ \frac{(1-\rho)^{\log n}}{\log n}, & h_{ji} = 1 \ \& \ v_i \in P_{j_s} \\ \frac{(1-\rho)^{h+1 \log n}}{\log n}, & 1 < h_{ji} = h \leq \sqrt{\frac{n}{\log n}} \ \& \ v_i \in P_{j_s} \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

Therefore, each node containing the content will serve only the nodes in the same cell with high probability, and the probability of being selected to serve the query initiated at the same cell is $\frac{1}{\rho \log n}$. Based on the bin-balls Theorem [17], the maximum number of queries served by a node will be $\frac{\log \log n}{\log \log \log n}$. Consequently, the maximum traffic load per source is $\gamma_{max} \frac{\log \log n}{\log \log \log n} = \frac{W \log \log n}{(1-\rho) \log n \log \log \log n}$. Therefore, to make sure that all the cells can support the stated throughput $\frac{1}{1-\rho} = 1 + \lambda/\mu$ is not allowed to exceed $\frac{\log n \log \log \log n}{\log \log n}$.

Finally each download of information will generate a traffic load on all the intermediate cells on the path from the source to the customer. However as stated in section III, the probability that the required content is discovered at distance of one hop is $1 - (1 - \rho)^{2 \log n}$ which is almost one for large n . So we may conclude that with high probability in sufficiently large networks no cell is working as relay or the number of transmissions passing through a cell as relay is close to zero.

Scenario i - Similar to scenario iii , the maximum traffic load is the load generated in a node when serving the requests. Here there are $n(1 - \rho)$ requests which will be served by $n\rho$ other nodes, so according to the bin-balls Theorem the maximum requests for a node will be in the order of $\frac{\log n}{\log \log n}$, which generates $\frac{W \rho \log n}{n(1-\rho) \log \log n} = \frac{W \lambda \log n}{n \mu \log \log n}$ traffic at the busiest node. This traffic does not exceed $\Theta(1)$ since the maximum requests to drop rate in this case is $\frac{n \log \log n}{\log n}$. ■

V. DISCUSSION AND FUTURE WORK

We studied the impact of caching with limited lifetime on the maximum capacity order in the grid and random networks where the received data is stored at the receivers and is shared with the other nodes as long as the node keeps the content. Figure 3 (a) shows the maximum throughput order for $\lambda/\mu = 7$ as a function of the network size. According to Theorem 1 and as can be observed from this figure, the maximum throughput capacity of the network in a grid network with the described characteristics is inversely proportional to the square root of the network size if the request rate and the cache timeout times are fixed. Similarly in the random network the maximum throughput is inversely proportional to the logarithm of the network size.

On the other hand with a fixed network size, if the ratio order of the request rate to the content time-out rate is greater than a threshold ($\Theta(n^{-1/2})$ in cases i, ii and $\Theta(\log^{-1} n)$ in case iii), most of the requests will be served by the caches and not the server, so increasing the request rate or decreasing the time-out rate will increase the probability of an intermediate cache having the content and reduces the number of hops needed to forward the content to the customer, and consequently increases the throughput (Figure 3 (b)). For request to time-out ratio orders less than these thresholds most of the requests are served by the main server (p_s approaches 1), so the maximum possible number of hops will be traveled by each content to reach the requester and the minimum throughput capacity ($\Theta(\frac{W}{n(1-\rho)})$ in cases i, ii , and $\Theta(\frac{W}{\sqrt{n \log n(1-\rho)}})$ in case iii) will be achieved.

Figures 4 (a),(b) respectively illustrate the total request rate and the total traffic generated in a fixed size network in scenario i for different request to time-out rate ratios. The total request rate in the network is the product of the number of requesting nodes and the rate at which each node is sending the request. The total traffic is the product of the total request rate and the number of hops between source and destination and the content size. Small λ/μ means that each node is sending requests with low rate, so fewer nodes have the content and more nodes are sending requests. In this case most of the requests are served by the server. The total request rate will increase by increasing the per node request rate. High λ/μ shows that each node is requesting the content with higher rate, so the number of cached content in the network is high and fewer nodes are requesting the content. Here most of the requests are served by the caches. The total request rate then is determined by the content drop rate. So for very large λ/μ , the

total request rate is the total number of nodes in the network times the drop rate ($n\mu$) and the total traffic is $n\mu B$.

However, when the network grows the traffic in the network will increase and the download rate will decrease. If we assume that the new requests are not issued in the middle of the previous download, the request rate will decrease with network growth. If the holding time of the contents in a cache increases accordingly the total traffic will not change, i.e. if by increasing the network size the requests are issued not as fast as before, and the contents are kept in the caches for longer times, the network will perform similarly.

Furthermore, if the ratio of request to content drop rates increases with network growth, higher throughput capacities may be achievable. For example in scenario *iii* if $\frac{\lambda}{\mu} = \Theta(\frac{\log n}{\log \log n})$, then the resulting throughput will be $\gamma_{max} = \Theta(\frac{W}{\log \log n}) \gg \Theta(\frac{W}{\log n})$. Note that according to Theorem 2, $\frac{\lambda}{\mu}$ is upper bounded by some values, so the achievable capacity will be upper bounded by $\Theta(\frac{W\sqrt{n} \log \log n}{\log n})$ (*i*) and $\Theta(W(\frac{\log \log n}{\log \log n} + \frac{1}{\log n}))$ (*iii*).

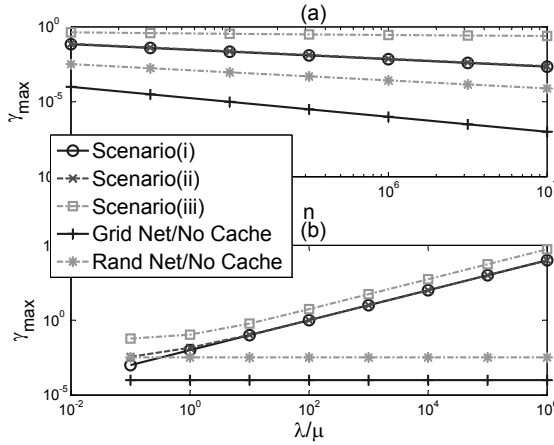


Fig. 3. Maximum download rate (γ_{max}) vs. (a) the number of nodes (n), (b) the request to time-out rate ratio (λ/μ).

In this work we have made several assumptions to simplify the analysis. For example, we assumed all the contents have the same characteristics (size, popularity). These assumptions will be relaxed in future work. We also assumed that the data is cached just in the receiver, and the requester downloads the data completely from one nearest content location. However, if all the intermediate nodes are allowed to cache the data and share it (as in the model of [18] for instance), or the node that needs the data can download each part of it from different nodes and makes a complete content out of the collected parts, higher capacities may be achievable. Proposing a caching and downloading scheme that can improve the capacity order is part of our future work.

REFERENCES

- [1] L. Zhang, D. Estrin, J. Bruke, V. Jacobson, J. Thornton, D. Smetters, B. Zhang, G. Tsudik, K. Claffy, D. Krioukov, D. Massey, C. Papadopoulos,

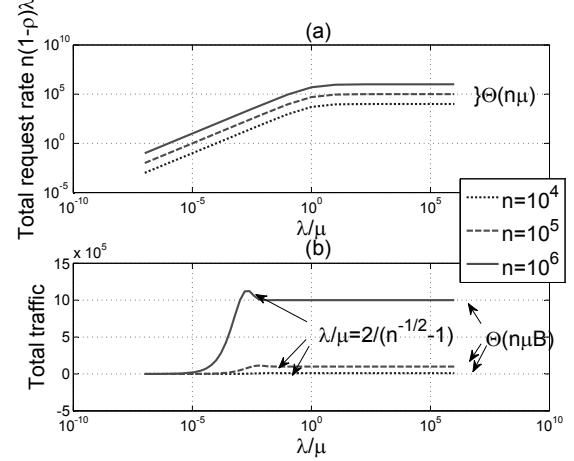


Fig. 4. (a) Total request rate in the network ($\lambda n(1-\rho)$), (b) Total traffic in the network ($B\lambda n(1-\rho)(p\bar{h} + p_s\bar{h}_s)$) vs. the request to time-out rate ratio (λ/μ).

- los, T. Abdelzaher, L. Wang, P. Crowley, and E. Yeh, "Named data networking (NDN) project," Oct. 2010.
- [2] "PURSUIT: Pursuing a pub/sub internet," <http://www.fp7-pursuit.eu/>, Sep. 2010.
- [3] B. Ahlgren, C. Dannewitz, C. Imbrenda, D. Kutscher, and B. Ohlman, "A survey of information-centric networking," *Communications Magazine, IEEE*, vol. 50, no. 7, July 2012.
- [4] V. Jacobson, D. K. Smetters, J. D. Thornton, M. F. Plass, N. H. Briggs, and R. L. Braynard, "Networking named content," in *ACM CoNEXT '09*. ACM, 2009, pp. 1–12.
- [5] H. Che, Z. Wang, and Y. Tung, "Analysis and design of hierarchical web caching systems," in *IEEE INFOCOM*, 2001, pp. 1416–1424.
- [6] E. Rosensweig, J. Kurose, and D. Towsley, "Approximate models for general cache networks," in *IEEE INFOCOM*, 2010, pp. 1–9.
- [7] E. J. Rosensweig and J. Kurose, "Breadcrumbs: Efficient, Best-Effort content location in cache networks," in *INFOCOM 2009, IEEE*. IEEE, Apr. 2009, pp. 2631–2635.
- [8] L. Yin and G. Cao, "Supporting cooperative caching in ad hoc networks," *IEEE Transactions on Mobile Computing*, no. 1, pp. 77–89, 2005.
- [9] M. Tortelli, I. Cianci, L. A. Grieco, G. Boggia, and P. Camarda, "A fairness analysis of content centric networks," Nov. 2011.
- [10] C. Westphal, "On maximizing the lifetime of distributed information in ad-hoc networks with individual constraints," in *MobiHoc '05*. New York, NY, USA: ACM, 2005, pp. 26–33.
- [11] H. Liu, Y. Zhang, and D. Raychaudhuri, "Performance evaluation of the cache-and-forward (CNF) network for mobile content delivery services," in *ICC Workshop*, 2009, pp. 1–5.
- [12] G. Carofiglio, M. Gallo, L. Muscariello, and D. Perino, "Modeling data transfer in content-centric networking," in *Teletraffic Congress (ITC), 2011 23rd International*. IEEE, pp. 111–118.
- [13] U. Niesen, D. Shah, and G. Wornell, "Caching in wireless networks," *IEEE Transactions on Information Theory*, 2011.
- [14] S. Gkitzenis, G. S. Paschos, and L. Tassioulas, "Asymptotic laws for content replication and delivery in wireless networks," in *INFOCOM, 2012 Proceedings IEEE*. IEEE, Mar., pp. 531–539.
- [15] F. Xue and P. Kumar, *Scaling Laws for Ad Hoc Wireless Networks: an Information Theoretic Approach*. Foundations and Trends in Networking, NOW Publishers, 2006.
- [16] M. D. Penrose, "The longest edge of the random minimal spanning tree," *The Annals of Applied Probability*, pp. 340–361, 1997.
- [17] M. Raab and A. Steger, "Balls into bins - a simple and tight analysis," in *Proceedings of the Second International Workshop on Randomization and Approximation Techniques in Computer Science*, 1998, pp. 159–170.
- [18] C. Westphal, K. Seada, C. Perkins, and R. Wakikawa, "An epidemiological study of information dissemination in mobile networks," in *IEEE SECON '09*, Jun. 2009.